

IN THE UNITED STATES PATENT AND TRADEMARK OFFICE

APPLICATION FOR LETTERS PATENT

Compressed File System For Non-Volatile RAM

Inventors:

Daniel J. Shoff

Jun Liu

&

John C. Southmayd

ATTORNEY'S DOCKET NO. MS1-430US

TECHNICAL FIELD

This invention relates generally to data storage mechanisms and, more particularly, to methods and arrangements for providing a compressed file system in non-volatile random access memory (NVRAM) or other like data storage mechanism.

BACKGROUND

During normal operations, information is stored within a computing device in the form of data. For example, in a typical personal computer (PC) data can be temporarily stored in a system or primary memory consisting of random access memory (RAM). Further data may be stored in a secondary memory, such as, a hard drive, a floppy drive, optical drive, etc. In such configurations, the secondary memory provides a non-volatile data storage alternative to the primary memory since the data stored in the primary memory would usually be lost if a power failure occurred. Solid-state non-volatile RAM (NVRAM), which is capable of retaining the stored data even after a power failure, is currently available but tends to be cost prohibitive for most PCs.

There are some computing devices, however, that can benefit from the use of NVRAM. One example is a set top box (STB) that is used to provide entertainment and related information services to subscribing users. Very little data storage is typically needed to provide such services, as compared to a full-fledged PC application. Thus, all or part of the primary memory in the STB can be NVRAM. This essentially eliminates the need to provide a secondary memory, such as, a hard drive, or the like. This can reduce the overall complexity, size and cost of the STB, for example.

1 Since NVRAM tends to be relatively expensive when compared to RAM
2 and/or conventional secondary memory, there is an attendant need to significantly
3 reduce or otherwise minimize the amount of NVRAM required within a given
4 computing device.
5

6 SUMMARY

7 Methods and arrangements are provided that significantly reduce or
8 otherwise minimize the amount of NVRAM required within a given computing
9 device.

10 For example, a novel data structure and management scheme are provided
11 in a manner that allows an NVRAM sector-based memory to appear as providing
12 significantly more storage space than it physically has. This is accomplished by
13 mapping a higher number of virtual sectors to a fewer number of physical sectors.
14 Data written to a plurality of virtual sectors is compressed and written to physical
15 sector(s). The information needed to associate the virtual and physical sectors can
16 be maintained in a virtual sector table within less expensive RAM. If power is lost
17 and the virtual sector table is no longer available in the RAM, then on power-up
18 the virtual sector table is recreated based in information that is imbedded within
19 the stored data structure in physical sectors of the NVRAM. The scheme
20 promotes data integrity by carefully controlling the compression and
21 decompression processes and providing data and operational step backup
22 information to insure that data within the NVRAM is not lost by a sudden power
23 loss, etc.
24
25

1 The above stated needs are met, for example, by a method that includes
2 receiving file system data from an operating system, storing the file system data in
3 a plurality of reserved sectors within a non-volatile memory, compressing the file
4 system data stored within in the plurality of reserved sectors to create a
5 compressed data block, and storing the compressed data block in at least one
6 physical subsector within the non-volatile memory. In this manner the operating
7 system can be provided with a plurality of operatively accessible virtual sectors
8 resulting in a virtual memory capacity that exceeds the actual physical capacity of
9 the non-volatile memory.

10 The method may further include, therefore, mapping the plurality of virtual
11 sectors to at least one physical subsector through a Virtual Sector Table (VST)
12 stored in a volatile memory and presenting the operating system with the VST. A
13 Sector Allocation Table (SAT) within the volatile memory can be used to map
14 physical subsectors to the VST. The SAT can be generated based at least on a
15 unique group identifier that is stored in each physical subsector associated with
16 storing the compressed data block. For example, the SAT can be generated during
17 a device initialization time.

18 The step of storing the compressed data block in at least one physical
19 subsector within the non-volatile memory may further include associating each
20 physical subsector with a unique group identifier. The step of storing the
21 compressed data block in at least one physical subsector within the non-volatile
22 memory may further include writing each physical subsector associated with the
23 compressed data block to the non-volatile memory in an a sequential order, but not
24 necessarily a contiguous order.
25

1 In still other implementations, the step of storing the compressed data block
2 in at least one physical subsector within the non-volatile memory further includes,
3 maintaining input/output (I/O) operation status information within the non-volatile
4 memory, during on-going I/O operations, that identifies the status of on-going I/O
5 operations with respect to data stored within non-volatile memory.

6 An arrangement for use in providing an application access to a non-volatile
7 memory is also provided. The arrangement includes an operating system and a
8 device driver. Here, the operating system is configured to exchange input/output
9 (I/O) requests with the application and exchange corresponding file system
10 requests with the device driver. The device driver is configured to store the file
11 system data received from the operating system in a plurality of reserved sectors
12 within the non-volatile memory, compress the file system data stored within in the
13 plurality of reserved sectors to create a compressed data block, and store the
14 compressed data block in at least one physical subsector within the non-volatile
15 memory. This arrangement can be included along with a processor configured to
16 run the operating system and the device driver, and a non-volatile memory that is
17 operatively coupled to the processor as part of a set top box or other like
18 computing device.

19 A non-volatile computer-readable medium having stored thereon a data
20 structure, is also provided. The data structure includes a raw sector map data and a
21 plurality of raw sectors as identified in the raw sector map that are configurable to
22 store uncompressed data during input/output (I/O) operations, I/O operation status
23 data identifying on-going I/O operation, and a plurality of subsectors contiguously
24 arranged and configured to support the identified on-going I/O operations by
25 storing compressed data blocks derived from the uncompressed data.

1
2
3 **BRIEF DESCRIPTION OF THE DRAWINGS**

4 A more complete understanding of the various methods and arrangements
5 of the present invention may be had by reference to the following detailed
6 description when taken in conjunction with the accompanying drawings wherein:

7 Fig. 1 is a block diagram depicting an exemplary computing device having
8 a non-volatile random access memory (NVRAM).

9 Fig. 2 is a block diagram depicting an exemplary functional arrangement of
10 a computing device as in Fig. 1.

11 Fig. 3 is a functional block diagram depicting an arrangement for
12 compressing file system data into an NVRAM as in Fig. 1.

13 Fig. 4 is an illustrative diagram depicting the layout of an NVRAM as in
14 Fig. 1.

15 Fig. 5 is an illustrative/functional block diagram file system data being
16 compressed and stored in an NVRAM as in Fig. 1.

17
18 **DETAILED DESCRIPTION**

19 The following description describes various methods and arrangements that
20 are implemented in a computing device running Microsoft® Windows® CE. This
21 is by way of example only. Those skilled in the art will clearly recognize that the
22 various methods and arrangements as described below and shown in the Drawings
23 can be used in any type of computing device having any type of operating system
24 that is configured to read/write data from/to a non-volatile memory arrangement.
25

1 Windows® CE is a scalable operating system (herein after, referred to as
2 the “ OS”) that is designed for a variety of embedded systems and products. Its
3 multithreaded, multitasking, fully preemptive operating environment is designed
4 specifically for hardware with limited resources, such as consumer electronic
5 devices, specialized industrial controllers, and embedded communications devices.
6 The OS also supports various hardware peripherals, devices, and networking
7 systems. These include keyboards, mouse devices, touch panels, serial ports,
8 Ethernet connections, modems, universal serial bus (USB) devices, audio devices,
9 parallel ports, printer devices, and storage devices, such as PC Cards.

10 With this in mind, Fig. 1 is a functional block diagram depicting an
11 exemplary computing device 20 having at least one processor 22 that is
12 operatively coupled to a memory 24 through at least one bus 26. Processor 22 is
13 further operatively coupled to at least one input/output device 28 through bus 26.
14 In this example memory 24 includes read only memory (ROM) 30, non-volatile
15 random access memory (NVRAM) 32, and random access memory (RAM) 34.
16 During normal operations, a sector allocation table (SAT) 36, a virtual sector table
17 (VST) 38 and a Group ID table 40 are generated and stored in RAM 34. These
18 tables and the information in them are described in greater detail in subsequent
19 sections. Those skilled in the art will recognize that computing device 20 may be
20 implemented through the use of one or more integrated circuits, etc.

21 Fig. 2 is a functional block diagram depicting an exemplary functional
22 arrangement 60 of computing device 20, for example, as in Fig. 1. Here, an
23 application program 62 is operatively configured to interface with an operating
24 system (OS) 64 and support reading and writing operations to a memory. OS 64 is
25

1 in turn operatively configured to interface with a device driver, in this case, a
2 block device driver, that further supports the read/write operation.

3 Block device drivers are provided for data storage devices that allow data to
4 be read or written only in blocks of a fixed size. Block devices do not allow
5 individual bytes of data to be read or written. Block sizes tend to be one or a few
6 kilobytes; some common sizes are 512 bytes, 1 KB, 2 KB, and 4 KB. Block
7 devices are ideally suited to mass storage and persistent storage applications, such
8 as disk and tape drives, or NVRAM.

9 As shown in Fig. 2, OS 64 includes a FAT File System (FATFS) 66 to
10 support an "NVFile" device driver 68. FATFS 66 essentially implements a logical
11 file system and provides an abstraction between files in the application name space
12 and devices in the device name space. For example, FATFS can exchange "Read
13 File ()" 70, "Write File ()" 72 and "IOCTL_DISK_FORMAT_VOLUME 74, as
14 applicable, with application 62. The non-volatile file ("NVFile") device driver 68,
15 in this example a block device driver, is responsible for guaranteeing safe I/O
16 operations even when device 20 is interrupted by a power cycle. NVFile 68
17 exchanges a read request "DISK_IOCTL_READ" 76 and write request
18 "DISK_IOCTL_WRITE" 78, as applicable, with FATFS 66. In the above
19 example, both FATFS 66 and NVFile 68 are implemented as data link lists (e.g.,
20 dll modules).

21 NVFile device driver 68 creates a storage device of NVRAM 32 to persist
22 raw data from the FATFS file system. As described herein, it is possible and
23 advantageous to extend the available storage space by compressing file system
24 data into NVRAM 32. This can be accomplished without making substantial
25 changes to the FATFS layer supplied by the core OS 64.

1 The successful implementation of a compressed file system as described
2 herein takes advantage of conventional software compression schemes, such as the
3 LZX library from the OS. This exemplary software compression scheme tends to
4 yields an average compression ratio of approximately 4:1.

5 The NVFile storage device in this example can be addressed by FAT-12
6 implementation using 512-byte sectors, and one sector per cluster, for example.
7 Here, FATFS 66 notifies its underlying block device driver NVFile 68 when
8 formatting or file (cluster) delete operations occur. FATFS 66 also performs a
9 low-level format whenever a media format operation is performed.
10 Communication between FATFS 66 and NVFile 68 insures correct reporting of
11 "disk full" status and the existence of free space. For example, FATFS 66 notifies
12 NVFile 68 when clusters are deleted, to insure that free space is reported
13 accurately and free clusters can be properly identified as free space.

14 Reference is now made to Fig. 3, which is a functional block diagram
15 depicting an arrangement for compressing file system data into NVRAM 32 as in
16 Fig. 1. As depicted, FATFS 66 provides a write request (78) that specifies a buffer
17 "pData" (e.g., a block of data) and a "Sector XXX". Sector XXX is then mapped
18 to a "virtual sector" within VST 38. Since the actual physical storage capacity of
19 NVRAM 32 is less than that supported by VST 38, a compression stage 80 is
20 introduced in NVFile to compress the "pData" and store a corresponding
21 compressed block of data on NVRAM 32.

22 This process saves NVRAM for data and is designed to avoid wasting
23 valuable NVRAM by persisting too much organizational information or tables of
24 pointers. Thus, in this example, tables and pointers are almost exclusively
25 maintained in RAM 34 rather than in NVRAM 32. Nevertheless, atomic/safe

1 operations are provided to maintain data and organizational integrity across
2 power-cycles. As will be appreciated, the methods and arrangements presented
3 herein also tend to avoid heap fragmentation.

4 Fig. 4 is an illustrative diagram depicting an exemplary layout of NVRAM
5 32. Here, NVRAM 32 includes a raw sector map 100, an I/O operation status 102,
6 a plurality of reserved raw sectors 104 (numbered 0-6), and a plurality of
7 subsectors 106 (numbered 0-N).

8 Let us assume, in this example, that NVRAM 32 is a 100-k byte memory
9 within a set top box. Here, the beginning of the 100-k block of NVRAM 32 starts
10 with a 28-byte raw sector map 100 (one DWORD for each reserved sector), then a
11 series of seven 512-byte raw uncompressed reserved sectors 104. Following the
12 end of the reserved sectors 104 are the 128-byte subsectors 106 to store
13 compressed blocks of data from the file system.

14 In raw sector map 100, either a virtual sector number or an "unused" value
15 is used to indicate what is stored in each reserved sector 104.

16 An I/O operation status 102 is provided between raw sector map 100 and
17 the first reserved sector 104. I/O operation status 102 stores the current state of
18 NVFile 68 in case of a power failure or normal power cycle during an I/O
19 operation. For example, I/O operation status 102 can include a flag indicating the
20 type of operation (write, delete, compress, etc.) and then some additional data
21 depending on the context/stage of the operation.

22 Since, as described below, compression will normally occur after every
23 fourth sector-write in this exemplary arrangement, seven reserved sectors 104 are
24 employed to cache data in the very likely event that a power cycle will occur while
25

1 the data are awaiting compression. Thus, all data from the file system is persisted
2 somewhere while it is being moved.

3 The seven reserved sectors 104 include additional space for use during
4 delete operations. For example, as part of a delete operation the compressed block
5 containing the deleted cluster must be decompressed and reassembled with a new
6 fourth sector. To provide for safe/atomic operation, all four sectors of that
7 particular compressed block are stored in reserved sectors 104 until such time as
8 the delete operation and recompression has been successfully completed.

9 The remaining space in NVRAM 32 is divided into 128-byte subsectors
10 where compressed data blocks are stored. Each 128-byte subsector will begin
11 with a group ID 108 (or marked "unused"), as shown in greater detail in
12 subsectors 106a and 106b. Group ID 108 will identify which subsectors 106
13 belong to which compressed block, since they will not necessarily be written
14 contiguously. In addition, as depicted in subsector 106a, the first subsector of a
15 compressed block within NVRAM 32 will have, following the group ID 108, four
16 virtual sector numbers 110 to identify where the data belongs in the real file
17 system.

18 Fig. 5 is an illustrative/functional block diagram file system data being
19 compressed and stored in subsectors 106. Here, it is assumed that write requests
20 #1, #2, and #3 have each been received and stored in raw sectors 0, 1, and 2,
21 respectively. When write request #4 is subsequently received, all four data blocks
22 are provided to a data compressor function 120 to produce a compressed block of
23 data. The compressed block of data is then provided to a splitter function that
24 stores pieces A, B and C of the compressed data in unused subsectors 106.
25

1 The Sector Allocation Table (SAT) 36 is essentially the "glue" relating the
2 compressed block subsectors 106 of NVRAM 32 to the Virtual Sector Table
3 (VST) 38 that is presented to the file system. Group IDs 108 are stored in
4 compressed subsectors to allow for generation of SAT 36 at boot-up/init time.
5 SAT 36 is basically a fixed-size mirror of the compressed-block subsectors 106 in
6 NVRAM 32. There is one SAT entry for each compressed subsector 106, and
7 each one will either be marked "unused" or indicate the next compressed subsector
8 106 in the chain. The last member of a chain will point to itself to mark the end.

9 SAT information is not be stored directly to NVRAM 32, but rather
10 inferred at boot time using Group ID 108 and the Virtual Sector Numbers 110 that
11 are stored in every physical compressed subsector 106. The Group ID number and
12 the sequential order in which subsectors having a common Group ID 108 are
13 written/read, will be used to generate the compressed subsector chain information
14 in SAT 36.

15 At boot-time, for example, each NVRAM subsector 106 is read, in order
16 from the beginning to the end of the subsector storage space. The first byte of
17 each subsector contains either an "unused" value, or a Group ID 108. For each
18 unique Group ID 108, every following subsector 106 with a matching Group ID
19 108 will be associated with that first subsector 106. The resulting chain of
20 associated subsectors 106, in the order found, will eventually comprise a
21 compressed block that can be reassembled and decompressed. For the first unique
22 Group ID found, an entry is made in the SAT that is later filled in with the offset
23 address of the next subsector 106 in that group. When the final subsector 106 is
24 found (determined only by reaching the end of the NVRAM and finding no more
25

subsectors in that group), the final SAT 36 entry is marked with its own address, indicating the end of the chain.

As described above, instead of mapping FATFS sectors directly to like-sized blocks of NVRAM 32, NVFile 68 presents VST 38 to FATFS while actually storing compressed blocks of data in the actual NVRAM subsectors 106. In order to minimize NVRAM overhead of this scheme, VST 38 is generated on the fly from this information at boot time, rather than stored in NVRAM in its entirety. In addition, a small portion of uncompressed data is stored in reserved raw sectors 104 of NVRAM 32 to persist sectors that have been written but not had a chance to compress yet.

To properly generate VST 38 at boot time, it is necessary to maintain and persist (during write operations) the relationship between the compressed data blocks and the virtual sectors. This relation is stored in the table of Group IDs 40, which map virtual sectors to groups of physical compressed subsectors 106.

When a new group of four sectors is compressed, they are assigned an unused Group ID 108. A table is stored at the beginning of the compressed data block, listing the virtual sector numbers represented by the compressed block. This relationship between the Group ID, physical subsectors, and virtual sectors numbers, is maintained at runtime. During initialization, Group ID table 40 is therefore dynamically generated, and maintained during all I/O operations.

By way of example, each entry in VST 38 can be a pointer to the first SAT entry corresponding to that virtual sector. For each compressed block containing four virtual sectors of data, there will be four VSTs pointing to its first SAT entry. The I/O operation can then start from that SAT entry and simply "walk" the pointers to retrieve and reassemble the entire compressed block.

1 The purpose of Group ID 108 is to store as little organizational (i.e., non-
2 data) information as possible in NVRAM. Every virtual sector in the VST must
3 point to a compressed block, but since each compressed block contains four virtual
4 sectors, otherwise unrelated virtual sectors will share the Group ID they reside in.

5 When four sectors are compressed, the new compressed block is assigned
6 the next-unused Group ID 108 from the Group ID table 40 in RAM 34. Group ID
7 108 itself is saved as the first byte in each physical sector that the compressed
8 block gets stored in.

9 During boot-up time, Group ID table 40 is re-created based on the stored
10 Group ID information. This can be accomplished by looping iteratively though
11 the NVRAM physical subsectors, matching up all of the sectors with others
12 sharing the same Group ID 108 to reestablish the associative relationship. Then,
13 the virtual sector numbers and Group Ids can be linked in the table in RAM 34.

14 Fig. 3 demonstrates how a sector write request from the file system is
15 routed through VST 38, and the data buffer is compressed and stored in a
16 compressed block of NVRAM 32.

17 As shown in the example of Fig. 5, to increase the efficiency of the data
18 compression, four sectors are compressed at a time rather than just one. However,
19 there is no guarantee that write operations will come in four-sector increments. As
20 a result, a buffer is persisted which can store uncompressed data that is waiting to
21 be compressed. The first "line" of physical NVRAM (raw sectors "0", "1", and
22 "2") is reserved for three raw sectors, and once those are filled they can be
23 combined with the fourth sector of incoming data, to produce a compressed block.

24 The compressed block will then be split so that it can fit into a number of
25 smaller physical memory subsectors 106. How many subsectors are required will

1 depend on the size of the final compressed block. Pieces of the compressed block
2 are then copied into NVRAM 32 and stored persistently.

3 The location of free physical memory sectors may be determined using an
4 iterative loop to find the next available unused subsector, if one is available. Since
5 it is difficult to know ahead of time the final size of a compressed block, it is
6 preferable to allow for fragmentation to occur so as to make better use of the
7 available memory space.

8 As part of the write operation, NVFile is also responsible for maintaining
9 the integrity of the Group ID table 40. As well as updating the dynamic Group ID
10 table 40 in RAM 34 to reflect the new write operation, it must also store the
11 correct Group ID 108 info in the subsectors 106 that are written.

12 After the compressed block is split into subsectors 106, they must be kept
13 in order so that the original compressed block can be recreated later for
14 decompression. From the beginning to the end of NVRAM 32, subsectors 106
15 belonging to any compressed block will be in order; there will be no actual
16 pointers stored in NVRAM 32 for purposes of maintaining subsector chains.

17 An incoming Read I/O operation will result in a hit on VST 38, where
18 actual file system sectors correspond to locations of NVRAM compressed blocks.
19 Since multiple virtual sectors are stored in a single compressed block, and the
20 compressed block must be decompressed to retrieve data, the entire compressed
21 block must be reassembled. Since this is true for any virtual sector residing within
22 the block, virtual sectors will all point to the beginning of their respective
23 compressed blocks. Further addressing is stored in the compressed block header
24 and decoded after decompression.

1 When a virtual sector hit point to any physical memory portion within a
2 compressed block, the first memory portion is read, and the chain followed to each
3 memory portion comprising the compressed block. The compressed block can
4 then be reassembled in a RAM memory buffer, decompressed, and the resulting
5 raw data returned to the file system.

6 There is a chance that a VST hit will point to an uncompressed raw sector
7 104 in the reserved line of physical NVRAM 32. In these cases, the uncompressed
8 data will simply be read straight through and returned to the file system.

9 When a virtual sector is deleted by FATFS, the compressed block in which
10 it resides is decompressed, and reassembled without the deleted data. If one or
11 more uncompressed sectors are already waiting in the reserved area 104, one of
12 them can be combined with the remaining three from the compressed block to
13 write a new block of four sectors. However, the existing compressed block cannot
14 be deleted until a new one is written. In order to insure safe/atomic operation, it is
15 important to reserve four additional uncompressed sectors 104 (i.e., raw sectors
16 "3", "4", "5", and "6") in the worst case but likely event that the first line of raw
17 sectors 104 (i.e., raw sectors "0", "1", and "2") are already full when a cluster is
18 deleted.

19 Deleted virtual sectors will be marked as unused in VST 38 and SAT 36,
20 and their former Group ID 108 is marked "unused".

21 A common scenario is to have a FATFS sector re-written with new data.
22 This operation takes place in much the same way as the aforementioned delete
23 operation. The main difference is that instead of finding a new virtual sector to
24 replace the deleted one, a new compressed block and subsector chain are stored
25 using the new data from the updated sector.

1 FATFS 66 may also be modified to notify NVFile 68 via an IOCTL when a
2 format operation takes place. This way NVFile can mark all of the clusters as
3 unused and reset all of the tables. It is important that this notification take place so
4 that the status of the used versus free sectors is consistent between the file system
5 and NVFile 68.

6 It should be understood that the methods and arrangements described herein
7 are not limited to data compression stages having compression ratios of 4:1,
8 memory layouts having seven reserved sectors, raw sectors of 512 bytes, or
9 compressed sectors of 128 bytes. For example, one or more of these parameters
10 can be higher or lower depending upon the computing device.

11 Thus, although some preferred embodiments of the various methods and
12 arrangements of the present invention have been illustrated in the accompanying
13 Drawings and described in the foregoing Detailed Description, it will be
14 understood that the invention is not limited to the exemplary embodiments
15 disclosed, but is capable of numerous rearrangements, modifications and
16 substitutions without departing from the spirit of the invention as set forth and
17 defined by the following claims.